

# Anomaly Detection Framework Based on Matching Pursuit for Network Security Enhancement

Rafał Renk, Witold Holubowicz  
ITTI Ltd., Poznań  
Adam Mickiewicz University, Poznań  
POLAND

[rafal.renk@itti.com.pl](mailto:rafal.renk@itti.com.pl) / [witold.holubowicz@itti.com.pl](mailto:witold.holubowicz@itti.com.pl)

## ABSTRACT

*In this paper, a framework for recognizing network traffic in order to detect anomalies is proposed. We propose to combine and correlate parameters from different layers in order to detect 0-day attacks and reduce False Positives. Moreover, we propose to combine statistical and signal-based features. The major contribution of this paper are: novel framework for network security based on the correlation approach as well as new signal based algorithm for intrusion detection using Matching Pursuit.*

## 1.0 INTRODUCTION AND MOTIVATION

Intrusion Detection Systems (IDS) are based on mathematical models, algorithms and architectural solutions proposed for correctly detecting inappropriate, incorrect or anomalous activity within a networked systems. Intrusion Detection Systems can be classified as belonging to two main groups depending on the detection technique employed: anomaly detection and signature-based detection. Anomaly detection techniques, that we focus on in our work, rely on the existence of a reliable characterization of what is normal and what is not, in a particular networking scenario. More precisely, anomaly detection techniques base their evaluations on a model of what is normal, and classify as anomalous all the events that fall outside such a model. If an anomalous behaviour is recognized, this does not necessarily imply that an attack activity has occurred: only few anomalies can be actually classified as attempts to compromise the security of the system.

Anomaly Detection Systems can be classified according to:

- the used algorithm,
- analyzed features of each packet singularly or of the whole connection,
- the kind of analyzed data - whether they focus on the packet headers or on the payload.

Most current IDS systems have problems in recognizing new attacks (0-day exploits) since they are based on the signature-based approach. In such mode, when system does not have an attack signature in database, such attack is not recognized. Another drawback of current IDS systems is that the used parameters and features do not contain all the necessary information about traffic and events in the network.

Therefore, in this paper we present the framework in which anomaly detection system based on correlation and diversity approaches are used, such as:

- item diversity - different network layers parameters are monitored and used. In such approach we do not have information from transport layer only - such information is merged/correlated with application layer events.

Report Documentation Page				Form Approved OMB No. 0704-0188		
Public reporting burden for the collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to a penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.						
1. REPORT DATE <b>NOV 2010</b>		2. REPORT TYPE <b>N/A</b>		3. DATES COVERED <b>-</b>		
4. TITLE AND SUBTITLE <b>Anomaly Detection Framework Based on Matching Pursuit for Network Security Enhancement</b>				5a. CONTRACT NUMBER		
				5b. GRANT NUMBER		
				5c. PROGRAM ELEMENT NUMBER		
6. AUTHOR(S)				5d. PROJECT NUMBER		
				5e. TASK NUMBER		
				5f. WORK UNIT NUMBER		
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) <b>ITTI Ltd., Poznan Adam Mickiewicz University, Poznan POLAND</b>				8. PERFORMING ORGANIZATION REPORT NUMBER		
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)				10. SPONSOR/MONITOR'S ACRONYM(S)		
				11. SPONSOR/MONITOR'S REPORT NUMBER(S)		
12. DISTRIBUTION/AVAILABILITY STATEMENT <b>Approved for public release, distribution unlimited</b>						
13. SUPPLEMENTARY NOTES <b>See also ADA564697. Information Assurance and Cyber Defence (Assurance de l'information et cyberdefense). RTO-MP-IST-091</b>						
14. ABSTRACT <b>In this paper, a framework for recognizing network traffic in order to detect anomalies is proposed. We propose to combine and correlate parameters from different layers in order to detect 0-day attacks and reduce False Positives. Moreover, we propose to combine statistical and signal-based features. The major contribution of this paper are: novel framework for network security based on the correlation approach as well as new signal based algorithm for intrusion detection using Matching Pursuit.</b>						
15. SUBJECT TERMS						
16. SECURITY CLASSIFICATION OF:				17. LIMITATION OF ABSTRACT <b>SAR</b>	18. NUMBER OF PAGES <b>10</b>	19a. NAME OF RESPONSIBLE PERSON
a. REPORT <b>unclassified</b>	b. ABSTRACT <b>unclassified</b>	c. THIS PAGE <b>unclassified</b>				

- correlation - correlation is used twofold (during decision):
- item both anomaly and signature-based approaches are correlated,
- parameters/features from various network layers are correlated,
- statistical and signal-based features are used and correlated.

## **2.0 TECHNICAL SOLUTION**

In this paper, a new solution for ADS system based on signal processing algorithm is presented. ADS analyzes traffic from internet connection in certain point of a computer network. The proposed ADS system uses redundant signal decomposition method based on Matching Pursuit algorithm. ADS based on Matching Pursuit uses Dictionary of Base Functions - BFD to decompose input 1D traffic signal (1D signal may represent packets per second) into set of based functions called also atoms. The proposed BFD has a ability to approximate traffic signal. Number and parameters of base functions was limited in order to shorten atom search time process

Since some attacks are visible only in specific layer (e.g. SQLIA), in our approach, we propose to use network parameters from different layers.

Transport layer, network layer and application layer parameters are used

In the further step, we use the presented parameters to calculate characteristics (features) of the observed traffic. Some of the parameters are used for statistical features calculation and/or for signal-based feature calculation respectively. Feature extraction methods are presented in the following subsections

### **2.1 Statistical Features**

The Chi-Square multivariate test for Anomaly Detection Systems can be represented by equation 1:

$$\chi^2 = \sum_{i=1}^p \frac{(X_i - \bar{X}_i)^2}{\bar{X}_i} \quad (1)$$

Where  $\mathbf{X} = (X_1, X_2, \dots, X_p)$  denote an observation of  $p$  variables from a process at time  $t$  and  $\bar{\mathbf{X}} = (\bar{X}_1, \bar{X}_2, \dots, \bar{X}_p)$  is the sample mean vector.

Using only the mean vector in Equation (1), cause that Chi-Square multivariate test detects only the mean shift on one or more of the variables.

### **2.2 Signal Processing Features**

Signal processing techniques have found application in Network Intrusion Detection Systems because of their ability to detect novel intrusions and attacks, which cannot be achieved by signature-based approaches. It has been shown that network traffic presents several relevant statistical properties when analyzed at different levels (e.g. self-similarity, long range dependence, entropy variations, etc.)

Approaches based on signal processing and on statistical analysis can be powerful in decomposing the signals related to network traffic, giving the ability to distinguish between trends, noise, and actual anomalous events. Wavelet-based approaches, maximum entropy estimation, principal component analysis techniques, and spectral analysis, are examples in this regard which have been investigated in the recent years by the research community. However, Discrete Wavelet Transform provides a large amount of coefficients which not necessarily reflect required features of the network signals.

Therefore, in this paper we propose another signal processing and decomposition method for anomaly/intrusion detection in networked systems. We developed original Anomaly Detection Type IDS algorithm based on Matching Pursuit.

In the rest of the paper, our original ADS method will be presented in details. Moreover, results of experimental setup will be given. We tested our method with standard traces in Worm detection scenario as well as in anomaly detection scenario. Discussion on redundant dictionary parameters and final conclusions will be provided.

## 2.2.1 Matching Pursuit

Matching Pursuit signal decomposition was proposed by Mallat and Zhang.

Matching Pursuit is a greedy algorithm that decomposes any signal into a linear expansion of waveforms which are taken from an over complete dictionary  $D$ . The dictionary  $D$  is an over complete set of base functions called also atoms.

$$D = \{\alpha_\gamma : \gamma \in \Gamma\} \quad (2)$$

where every atom  $\alpha_\gamma$  from dictionary has norm equal to 1:

$$\|\alpha_\gamma\| = 1$$

$\Gamma$  represents set of indexes for atom transformation parameters such as translation, rotation and scaling.

Signal  $s$  has various representations for dictionary  $D$  Signal can be approximated by set of atoms  $\alpha_k$  from dictionary and projection coefficients  $c_k$  :

$$s = \sum_{n=0}^{|D|-1} c_k \alpha_k \quad (3)$$

To achieve best sparse decomposition of signal  $s$  (min) we have to find vector  $c_k$  with minimal norm but sufficient for proper signal reconstruction. Matching Pursuit is a greedy algorithm that iteratively approximates signal to achieve good sparse signal decomposition. Matching Pursuit finds set of atoms  $\alpha_{\gamma_k}$  such that projection of coefficients is maximal. At first step, residual  $R$  is equal to the entire signal  $R_0 = s$ .

$$R_0 = \langle \alpha_{\gamma_0}, R_0 \rangle \alpha_{\gamma_0} + R_1 \quad (4)$$

If we want to minimize energy of residual  $R_1$  we have to maximize the projection.  $\langle \alpha_{\gamma_0}, R_0 \rangle$  At next step we must apply the same procedure to  $R_1$ .

$$R_1 = \langle \alpha_{\gamma_1}, R_1 \rangle \alpha_{\gamma_1} + R_2 \quad (5)$$

Residual of signal at step  $n$  can be written as follows:

$$R^n s = R^{n-1} s - \left\langle R^{n-1} s \mid \alpha_{\gamma_k} \right\rangle \alpha_{\gamma_k} \quad (6)$$

Signal  $s$  is decomposed by set of atoms:

$$s = \sum_{k=0}^{N-1} \left\langle \alpha_{\gamma_k} \mid R^n s \right\rangle \alpha_{\gamma_k} + R^n s \quad (7)$$

Algorithm stops when residual  $R^n s$  of signal is lower then acceptable limit.

### 2.2.2.1 Our Approach to Intrusion Detection Algorithm

In basic Matching Pursuit algorithm atoms are selected in every step from entire dictionary which has flat structure. In this case algorithm causes significant processor burden. In our coder dictionary with internal structure was used.

Dictionary is built from:

- Atoms
- Centered Atoms

Centered atoms groups such atoms from  $D$  that are as more correlated as possible to each other. To calculate measure of correlation between atoms function  $o(a, b)$  can be used [2].

$$o(a, b) = \sqrt{1 - \left( \frac{|\langle a, b \rangle|}{\|a\|_2 \|b\|_2} \right)^2} \quad (8)$$

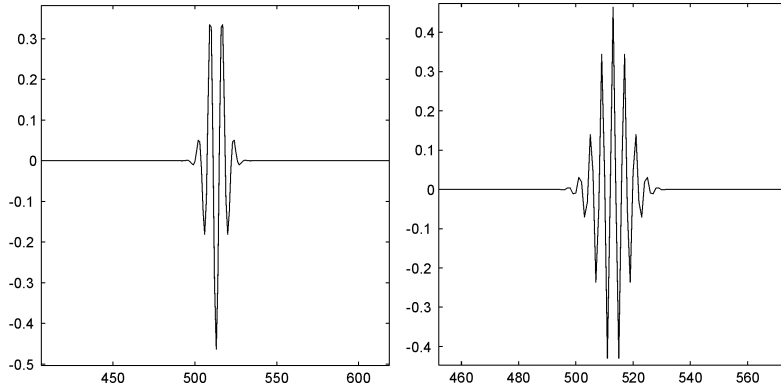
The quality of centered atom can be estimated according to :

$$O_{k,l} = \frac{1}{|LP_{k,l}|} \sum_{i \in LP_{k,l}} o(A_{c(i)}, W_{c(k,l)}) \quad (9)$$

$LP_{k,l}$  is a list of atoms grouped by centered atom.  $O_{k,l}$  is mean of local distances from centered atom  $W_{c(k,l)}$  to the atoms  $A_{c(i)}$  which are strongly correlated with  $A_{c(i)}$ .

Centroid  $W_{c(k,l)}$  represents atoms  $A_{c(i)}$  which belongs to the set  $i \in LP_{k,l}$ . List of atoms  $LP_{k,l}$  should be selected according to the Equation :

$$\max_{i \in LP_{k,l}} o(A_{c(i)}, W_{c(k,l)}) \leq \min_{t \in D \setminus LP_{k,l}} o(A_{c(t)}, W_{c(k,l)}) \quad (10)$$



**Figure 1: Example atom from dictionary**

In the proposed IDS solution 1D real Gabor base function (Equation was used to build dictionary (11))

$$\alpha_{u,s,\xi,\phi}(t) = c_{u,s,\xi,\phi} \alpha\left(\frac{t-u}{s}\right) \cos(2\pi\xi(t-u) + \phi) \quad (11)$$

where:

$$\alpha(t) = \frac{1}{\sqrt{s}} e^{-\pi^2 t^2} \quad (12)$$

$c_{u,s,\xi,\phi}$  - is a normalizing constant used to achieve atom unit energy,

In order to create over complete set of 1D base functions dictionary  $D$  was built by varying subsequent atom parameters: Frequency  $\xi$  and phase  $\phi$ , Position  $u$ , Scale  $s$ .

Base functions dictionary  $D$  was created with using 10 different scales (dyadic scales) and 50 different frequencies. In **Error! Reference source not found.** example atoms from dictionary  $D$  are presented.

## 2.2 Experimental Results

Percentage of the recognized anomalies as a function of encoded atoms from Dictionary of Base Functions is presented in Figure 2. Five dictionaries with different parameters (different number of scales and frequencies) were used in our ADS system.

Percentage of the recognized anomalies for Dictionary of Base Functions with approximately constant number of atoms is presented in Figure 3. In this case we try to leave approximately constant number of atoms in dictionary but with different proportions of scales and frequencies.

**Table 1: Matching Pursuit Mean Projection for TCP trace. Traces are analysed with the use of 20 minutes windows**

TCP trace	Window1 MP-MP	Window2 MP-MP	Window3 MP-MP	Mean. MP-MP for trace	Mean MP- MP for normal trace
Mawi 2004.03.06 tcp	210,34	172,58	239,41	245,01	240,00
Mawi 2004.03.13 tcp	280,01	214,01	215,46	236,33	240,00
Mawi 20.03.2004 (attacked: worm Witty)	322,56	365,24	351,66	346,48	240,00
Mawi 25.03.2004 (attacked: worm Slammer)	329,17	485,34	385,50	400,00	240,00

**Table 2: Matching Pursuit Mean Projection for UDP trace. Traces are analysed with the use of 20 minutes windows**

UDP trace	Window1 MP-MP	Window2 MP-MP	Window3 MP-MP	Mean. MP-MP for trace	Mean MP- MP for normal trace
Mawi 2004.03.06 tcp	16,06	13,80	17,11	15,65	16,94
Mawi 2004.03.13 tcp	20,28	17,04	17,40	18,24	16,94
Mawi 20.03.2004 (attacked: worm Witty)	38,12	75,43	61,78	58,44	16,94
Mawi 25.03.2004 (attacked: worm Slammer)	56,13	51,75	38,93	48,93	16,94

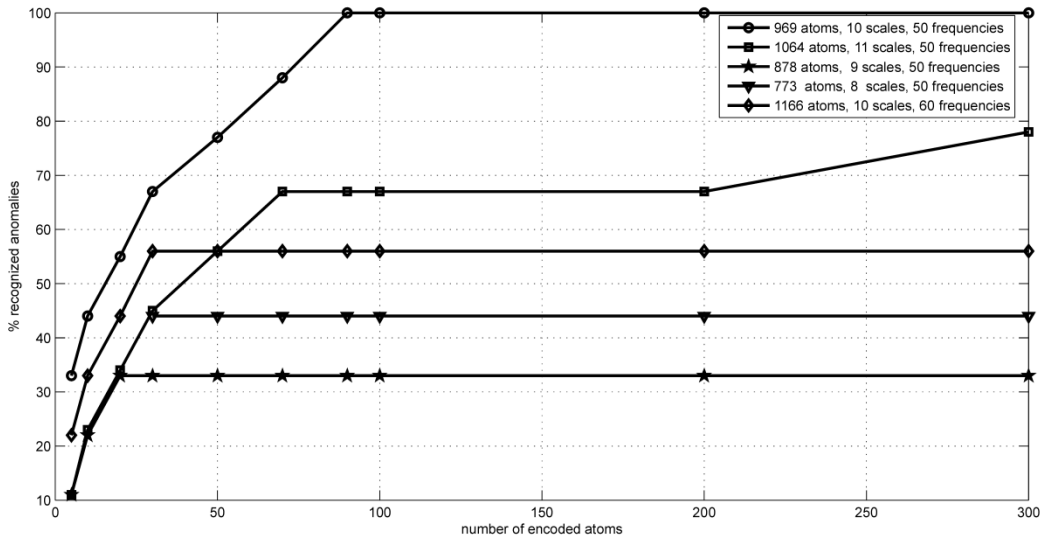
**Table 3: Matching Pursuit Mean Projection for TCP trace (traces consist of DDoS SYN Flood attacks). Traces are analysed with the use of 20 minutes windows**

TCP trace	Window1 MP-MP	Window2 MP-MP	Window3 MP-MP	Mean MP- MP for trace	Mean MP-MP for normal trace
One hour trace from unina1 [17]	1211	3271	3007	2496,333	860,00
One hour trace from unina2 [17]	1906	1804	1251	1653,667	860,00

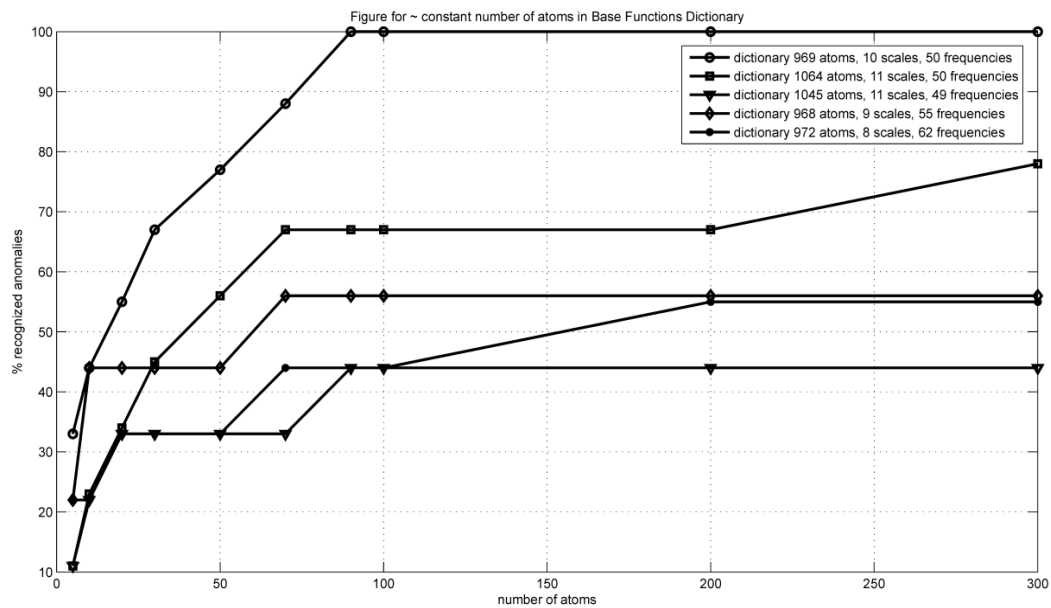
In Table 1,2,3,4 there are example results taken from our ADS system. Traffic traces were analysed by proposed ADS with the use of 20 minutes windows (most attacks (more than 80%) last no longer than 20 minutes). In every window we calculate Matching Pursuit Mean projection parameter in order to recognize suspicious traffic behaviour. Analysed traces are infected by worms (Table1,2), DDos(Table 4) and DDoS SYN Flood (Table 3) attacks.

**Table 4: Matching Pursuit Mean Projection for TCP trace (traces consist of DDoS attacks).  
Traces are analysed with the use of 20 minutes windows**

TCP trace	Window1 MP-MP	Window2 MP-MP	Window3 MP-MP	Mean. MP- MP for trace	Mean MP- MP for normal trace
Backscatter 2008.11.15	147,64	411,78	356,65	305,35	153,66
Backscatter 2008.08.20	208,40	161,28	153,47	174,38	153,66



**Figure 2: Percentage of the recognized anomalies as a function of encoded atoms**



**Figure 3: Percentage of the recognized anomalies for Dictionary of Base Functions with approximately constant number of atoms**



### **3.0 CONCLUSIONS**

In this paper a framework for recognizing attacks and anomalies in the computer networks is presented. Our methodology is based on both statistical and signal based features. The major contribution and innovation is the application of Matching Pursuit algorithm to calculate network traffic features. The effectiveness of the proposed approach has been proved in attack and anomaly detection scenarios. Our framework can be applied to enhance military networks since it uses signal-based features. Such features can be calculated for encrypted traffic since flow characteristics are extracted without considering the payload. Future work focuses on algorithms optimization so that our framework can be applied to real-time network security enhancement.

### **4.0 ACKNOWLEDGMENT**

The research leading to these results has received funding from the European Community's Seventh Framework Programme (FP7/2007-2013) under grant agreement no. 216585 (INTERSECTION Project).

### **BIBLIOGRAPHY**

1. Esposito M., Mazzariello C., Oliviero F., Romano S.P., Sansone C., Real Time Detection of Novel Attacks by Means of Data Mining Techniques. ICEIS (3) 2005: 120-127.
2. Esposito M., Mazzariello C., Oliviero F., Romano S.P., Sansone C., Evaluating Pattern Recognition Techniques in Intrusion Detection Systems. PRIS 2005: 144-153.
3. FP7 INTERSECTION Project, Deliverable D.2.1: SOLUTIONS FOR SECURING HETEROGENEOUS NETWORKS: A STATE OF THE ART ANALYSIS.
4. FP7 INTERSECTION (INfrastructure for heTERogeneous, Resilient, Secure, Complex, Tightly Inter-Operating Networks) Project Description of Work.
5. C.-M. Cheng, H. T. Kung, K.-S. Tan, Use of spectral analysis in defense against DoS attacks, IEEE GLOBECOM 2002, pp. 2143-2148.
6. P. Barford, J. Kline, D. Plonka, A. Ron, A signal analysis of network track anomalies, ACM SIGCOMM Internet Measurement Workshop 2002.
7. P. Huang, A. Feldmann, W. Willinger, A non-intrusive, wavelet-based approach to detecting network performance problems, ACM SIGCOMM Internet Measurement Workshop, Nov. 2001.
8. L. Li, G. Lee, DDos attack detection and wavelets, IEEE ICCCN03, Oct. 2003, pp. 421-427.
9. A. Dainotti, A. Pescapé, G. Ventre, Wavelet-based Detection of DoS Attacks, 2006 IEEE GLOBECOM - Nov 2006, San Francisco (CA, USA).
10. S. Mallat and Zhang Matching Pursuit with time-frequency dictionaries. IEEE Transactions on Signal Processing., vol. 41, no 12, pp. 3397-3415, Dec 1993.
11. J.A. Troop. Greed is Good: Algorithmic Results for Sparse Approximation. IEEE Transactions on Information Theory., vol. 50, no. 10, October 2004
12. R. Gribonval Fast Matching Pursuit with a Multiscale Dictionary of Gaussian Chirps. IEEE Transactions on Signal Processing., vol. 49, no. 5, may 2001.

13. P. Jost, P. Vandergheynst and P. Frossard Tree-Based Pursuit: Algorithm and Properties. Swiss Federal Institute of Technology Lausanne (EPFL), Signal Processing Institute Technical Report., TR-ITS-2005.013, May 17th, 2005.
14. A. Dainotti, A. Pescapé, G. Ventre, Worm Trac Analysis and Characterization, Proceedings of ICC, IEEE CS Press, 1435-1442, 2007.
15. WIDE Project: MAWI Working Group Traffic Archive at [tracer.csl.sony.co.jp/mawi/](http://tracer.csl.sony.co.jp/mawi/)
16. The CAIDA Dataset on the Witty Worm - March 19-24, 2004, Colleen Shanon and David Moore, [www.caida.org/passive/witty](http://www.caida.org/passive/witty).
17. Universita' degli Studi di Napoli "Federico II" (Italy), Network Tools and traffic traces, <http://www.grid.unina.it/Traffic/Traces/ttraces.php>

